



# Real-Time Video Quality Assessment in Packet Networks: A Neural Network Model

Samir Mohamed, Gerardo Rubino, Francisco Cervantes, Hossam Afifi

## ► To cite this version:

Samir Mohamed, Gerardo Rubino, Francisco Cervantes, Hossam Afifi. Real-Time Video Quality Assessment in Packet Networks: A Neural Network Model. [Research Report] RR-4186, INRIA. 2001. inria-00072437

**HAL Id: inria-00072437**

**<https://inria.hal.science/inria-00072437>**

Submitted on 24 May 2006

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

***Real-time Video Quality Assessment in Packet  
Networks:  
A Neural Network Model***

Samir Mohamed Gerardo Rubino

Francisco Cervantes Hossam Afifi

**N°4186**

Mai 2001

\_\_\_\_\_ THÈME 1 \_\_\_\_\_



***apport  
de recherche***





## Real-time Video Quality Assessment in Packet Networks: A Neural Network Model

Samir Mohamed \* Gerardo Rubino †  
Francisco Cervantes ‡ Hossam Afifi §

Thème 1 — Réseaux et systèmes  
Projet ARMOR

Rapport de recherche n° 4186 — Mai 2001 — 21 pages

**Abstract:** There is a great demand to assess video quality transmitted in real time over packet networks, and to make this assessment in real time too. Quality assessment is achieved using two types of methods: objective or subjective. Subjective methods give more reliable results than objective methods; the latter do not always correlate well with human perception. Unfortunately, subjective methods are not suitable to real-time applications and are very difficult to carry out. In this paper, we show how Artificial Neural Networks (ANN) can be used to mimic the way by which a group of human subjects assess video quality when this video is distorted by certain quality-affecting parameters (e.g. packet loss rate, loss distribution, bit rate, frame rate, encoded frame type, etc.). Our method can be used to measure in real time the subjective video quality with very good precision. In order to illustrate its applicability, we chose to assess the quality of video flows transmitted over IP networks and we carried out subjective quality tests for video distorted by variations of those parameters.

**Key-words:** Packet video, Neural Networks, Real-time video transmission, Video quality assessment, Video signal characterization.

(Résumé : *tsvp*)

This work was partially supported by the European ITEA Project 99011 “RTIPA” (Real-Time Internet Platform Architectures).

\* IRISA/INRIA, Campus du Beaulieu, 35042 Rennes, France. Email:Samir.Mohamed@irisa.fr

† IRISA/INRIA, Campus du Beaulieu, 35042 Rennes, France. Email:Gerardo.Rubino@irisa.fr

‡ Instituto Tecnológico Autónomo de México (ITAM), Mexico. Email:cervante@itam.mx

§ INT Evry, rue Charles Fourier, 91011 Evry, France. Email:Hossam.Afifi@int-evry.fr

## Évaluation en temps réel de la qualité vidéo dans les réseaux de paquets : une approche avec des réseaux de neurones

**Résumé :** Il y a une grande demande pour évaluer la qualité de la vidéo transmise en temps réel au-dessus des réseaux de paquets, et pour faire cette évaluation également en temps réel. Ceci est réalisé en utilisant deux types de méthodes : les méthodes objectives ou les méthodes subjectives. Les méthodes subjectives donnent des résultats plus fiables que les méthodes objectives ; ces derniers ne s'ajustent toujours pas bien avec la perception humaine. Malheureusement, les méthodes subjectives ne sont pas appropriées aux applications temps réel et sont très difficiles à mettre en œuvre. Dans cet article, nous montrons comment les réseaux de neurones artificiels (ANN) peuvent être employés pour imiter la façon par laquelle un groupe de sujets humains évalue la qualité de la vidéo déformée par certains paramètres (par exemple, le taux de perte de paquet, la distribution de perte, le débit du flux, la cadence de trame, etc.). Notre méthode peut être employée pour mesurer en temps réel la qualité subjective de la vidéo avec une très bonne précision. Afin d'illustrer son applicabilité, nous avons choisi d'évaluer la qualité des séquences vidéo transmises au-dessus des réseaux IP et nous avons effectué des tests subjectifs de qualité pour la vidéo déformée à cause de variations des valeurs de ces paramètres.

**Mots-clé :** Paquets vidéo, réseaux de neurones, transmission de vidéo temps réel, évaluation de la qualité de la vidéo, réseaux à commutation de paquets.

## 1 Introduction

Many real-time video transmission applications over the Internet have appeared in the last few years. We find today: video phones, video conferencing, video streaming, tele-medical applications, distance learning, telepresence, video on demand, etc., with a different number of requirements in bandwidth and perceived quality. This gives rise to a need for assessing the quality of the transmitted video in real time. Traditionally, the quality is measured as a function of the encoding algorithms, without taking into account the effect of packet networks' parameters.

Quality in encoders is also a way to "fit" the stream into the available global channel bandwidth. In video codecs using temporal compression (ex: MPEG, H.261, H.263), a quality factor parameter is usually used to reduce the output stream bandwidth and to reduce, in the same time, the assessed quality (yet before any transmission). An overview of transferring real-time video on packet networks as well as the most used video codecs is given in [3]. Now, if we consider the parameters that affect the quality (quality-affecting parameters) of video transmission over packet networks, we can classify them as follows:

- Coding and compression parameters: They control the amount of quality losses that happen during the encoding process; so they depend on the type of the encoding algorithm (MPEG, H26x, etc.), the output bit rate, the frame rate (the number of frames per sec.), the temporal relation among frame kinds (I, B, P, etc. frames), etc. [10], [15], [28];
- Network parameters: They result from packetization of the video stream [1] and the transmission in real-time, such as the packet loss rate, the loss distribution, the delay, the delay variation (jitter), etc. [11], [4], [6].
- Other parameters like the nature of the scene (e.g. amount of motion, color, contrast, image size, etc.) [10], [27], [26].

Since in this paper we concentrate only on pure video applications, we do not take into account parameters such as lip synchronization or other audio aspects. It is clear that quality is not linearly proportional to the variation of these parameters. The determination of the quality is a very complex problem, and there is no mathematical model that can take into account the effects of all these parameters.

There are two approaches to measure the quality: either by objective methods or subjective methods. The objective methods [29] measure quality based on mathematical analysis that compare original and distorted video sequence. Some existing methods are MSE (Mean Square Error) or PSNR (Peak Signal to Noise Ratio) which measures the quality by a simple difference between frames. There are other methods that are much more complicated like the moving picture quality metric (MPQM), and the normalized video fidelity metric (NVFM) [28].

On the other hand, subjective quality assessment methods [22] measure the overall perceived video quality. They are carried out by human subjects. The most commonly used for video quality evaluation is the Mean Opinion Score (MOS), recommended by the ITU. It consists in having  $n$  subjects viewing the distorted video sequences in order to rate their quality, according to a predefined quality scale. That is, human subjects are trained to “build” a mapping between the quality scale and a set of processed video sequences.

Although MOS studies have served as the basis for analyzing many aspects of signal processing, they present several limitations: a) very stringent environments are required; b) the process can not be automated; c) it is very costly and time consuming to repeat it frequently. Consequently, it is impossible to use it in real-time quality assessment.

On the other hand, the disadvantages of objective methods are: a) they do not correlate well with human visual perception ; b) they require high calculation power, and are time consuming (they usually operate at the pixel level); c) it is very hard to adapt them to real-time quality assessment, as they work on both the original video sequence and the transmitted/distorted one; d) as stated before, it is difficult to build a model that takes into account the effect of many quality-affecting parameters, specially network parameters. Then, instead of looking for algorithms to objectively measure video quality, why do not we build a hybrid system that takes into consideration subjective measurements, and which behavior is close to that of performing the same task? In this paper, we address this question by describing a method for developing such an automaton.

Our problem has two aspects: first, a classification one, to map the non-linear relation between the parameters and the quality; second, a prediction one, to evaluate the quality as a function of the quality-affecting parameters in an operational environment. We believe that artificial neural networks (ANN) are an appropriate tool to solve this two-fold problem [23], [18]. We illustrate our approach by building a system that takes advantage of the benefits offered by ANN to capture the nonlinear mapping between several non-subjective measures (i.e. the quality-affecting parameters) of video sequences transmitted over packet switched networks and the quality scale carried out by a group of humans subjects during an MOS experiment.

The structure of this paper is as follows. The next Section situates our study in its context, by describing related works. Section C presents our proposal in detail. Section D is dedicated to an overview of subjective quality and the computation of MOS. ANN are briefly described in Section E. Our results are the object of Section F. The last Section presents our conclusions and future research directions.

## 2 Related Works

In previous work, we showed how to use ANN to measure in real-time audio quality when this audio is transmitted over packet networks [20]. Based on this technique, we developed a new control mechanism that permits a better use of bandwidth and the delivery of the best possible audio quality given the current network situation [19].

The work in [21] presents a methodology for video quality assessment using objective parameters based on image segmentation. An image encoded by MPEG-2 is segmented into three regions: plane, edges, and texture; then, a set of objective parameters is assigned to each region. After that, a perception-based model is defined by computing the relationship between objective measures and results of subjective tests. This technique suffers from the drawbacks described in the previous Section, which are in fact common to most objective methods.

In [6], the authors study the effect of both loss and jitter on the perceptual quality of video. They argue that, if there is not a mechanism to mask the effect of jitter, the perceived quality degrades in the same way as it degrades with losses. While in [11], [27], [26] and [14], the authors analyze the effect of audio synchronization on the perceived video quality; they quantify the benefits of audio synchronization on the overall quality of the flow.

The main goal of [10] is to study the effect of the frame rate for different standard video sequences on the overall perceived quality. The work presented in [4] is mainly a study of the packet loss effects on MPEG video streams. The authors show also the effect of loss rate on the different types of MPEG frames. While in [28], a study of the effect of the bit rate on the objective quality metrics (PSNR, NVFM, and MPQM) is presented. The effect of the number of consecutively lost packets on the video quality is analyzed in [11]. In addition, the authors of [15] study the effect of packet size and the distribution of I-frames in the layered video transmission over IP networks. In [5] the analysis goes deeper: the authors present a study of the effect of motion on the perceived video quality.

In [16], the authors present how to use ANN to predict packet loss during real-time video transmission over packet networks as a function of the inter-packet delay variation. ANN are used also in video compression with compression ratio that goes from 500:1 to 1000:1 for moving gray-scale images and full-color video sequences respectively [7]. They are also used as decoders for error correcting codes in noisy communication channels, which reduces the error probability to zero [2]. Furthermore, they are used in a variety of image processing techniques which go from image enlargement and fusion to image segmentation [9].



### 3 Description of our Method

In this Section, we describe the overall steps that should be followed in order to build a tool to automatically assess in real time the subjective quality of real-time video transmitted over packet networks. The aim of this method is to use ANN to model and evaluate in real time how human subjects estimate video quality when distorted by changes in the quality-affecting parameters. In other words, such a tool will emulate a subjective MOS test carried out by a group of human subjects.

The overall procedure is summarized in Figure 1. We start by defining a set of static information that will affect the general quality perception. We must choose the most effective quality-affecting parameters corresponding to the type of video application and to the network supporting the transmission (see Section F.2).

Once the quality-affecting parameters are identified, for each one we should find the two extremes and the most frequent occurrences of its values. This can be done either by real measurement or using simulation techniques. For example, if the percentage loss rate is expected to vary from 0 to 10 %, then we may use 0, 1, 2, 5, and 10 % as the typical values for the loss rate parameter. If we call configuration of the set of quality-affecting parameters a set of values for each one, the total number of possible configurations is usually large. We must then select a part of this large cardinality set, which will be used as the input data of the ANN in the learning phase.

Depending on the transmission configuration, a simulation environment or a testbed should be implemented. This environment is used to send video sequences from the source to the destination and to control the underlying packet network. Every configuration in the defined input data must be mapped into the system composed of the network, the source and the receiver. For example, in our experiments with IP networks, the source controls bit rate, the frame rate and the encoding algorithm, and it sends RTP video packets; the router controls the loss rate, the loss distribution, the delay and the jitter; the destination stores the transmitted video sequence and collects the corresponding values of the parameters. Of course, one can generate the distorted signal by simulation. Then, by operating the testbed or the artificial simulation, we produce and store a set of distorted signals (this is the Video Database in Figure 1), along with their corresponding values of the parameters.

After completing the video database, a subjective quality test should be carried out. There are several subjective quality methods in the recommendations of the ITU-R [12] (see Section D.1). The most suitable one is Degradation Category Rating (DCR), discussed in Section D. The video database should be shuffled in such a way to avoid the effect of the last sequence on the judgment of the current one. A group of  $n$  human subjects is then invited to evaluate the quality of the distorted video sequences (i.e. every subject gives each video sequence a score

from the predefined quality scale). The subjects should not establish any relation between the sequences and the corresponding parameters' values.

The next step is to calculate the MOS values for all the video sequences. Based on the results obtained by the human subjects, a prescreening and statistical analysis may be carried out to remove the grading of the individuals suspected to give unreliable results. Details about subjective quality methods and MOS calculation are given in Section D. After that, we store the MOS values and the corresponding parameters' values in another database (Quality Database).

After that, a suitable neural network architecture and a training algorithm should be selected. We chose, as in other applications fields, a three-layered feedforward network and the backpropagation training. The training database can be divided into two parts: one to train the ANN and the other to test its accuracy. The trained ANN will emulate the subjective quality measure for any given values of the parameters (not necessary among the training database). The overall procedure should be repeated, if necessary, to improve the ANN's accuracy in evaluating video quality.

Once a stable neural network configuration is obtained, the ANN's architecture and weights can be extracted in order to build a concise tool. We decompose such a tool into two parts. The first one collects the values of the quality-affecting parameters. The second part is the trained ANN that will take the given values of the chosen quality-affecting parameters and correspondingly computes the subjective MOS quality score.

All the above steps are summarized into the four parts shown in Figure 1. In the first part, we have to identify the quality-affecting parameters, their ranges and values, to choose the original video sequences, and to produce the distorted video database. In the second part, the subjective quality test is carried out for the distorted video database and the MOS scores (together with their statistical analysis) are calculated in order to form the video quality database. While in the third part, we select the ANN architecture and learning algorithm, and train and test it using the video quality database. Finally, in the fourth step, we implement the final tool that consists of the two parts (parameters collection and MOS evaluation), in order to obtain a real-time quality evaluator.

### 3.1 Operating Mode

Real-time video applications can be considered one-way sessions (i.e. they consist of a sender that produces the video and a receiver that consumes it). This behavior is different from that of audio applications. Indeed, in audio, the interactivity may produce some other parameters (e.g. echo, crosstalk effect, number of participating sources, etc.) that affect the overall quality [20].

In the operating mode when integrated in a video system, our tool will act as shown in Figure 2. In the sender the video source is encoded and affected by some parameters. Then it is packetized and sent by the transport protocol (e.g. RTP/RTCP) to the receiver. Here again the

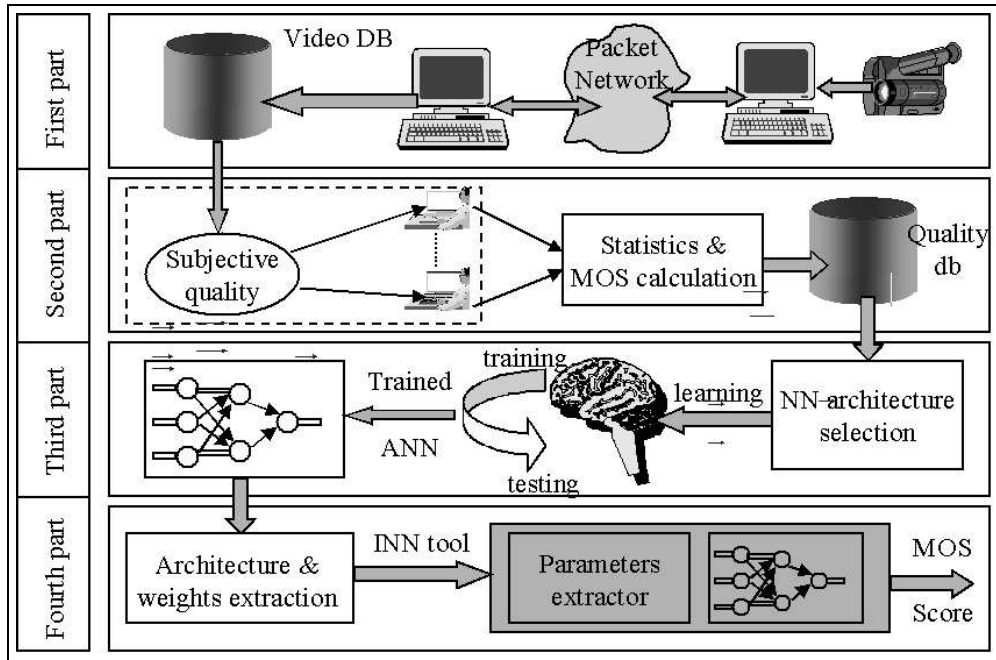


Figure 1: The overall architecture of the proposed method to evaluate real-time video quality in real time.

video quality may be degraded by certain parameters. In the receiver, the flow is decoded and displayed to the end-user.

The interaction between our tool and the other elements is as follows. The parameters' collector part probes all the working parameters from the encoder, decoder, packet network and the transport protocol. Then the trained ANN part evaluates video quality as a function of these parameters. In this way, the end-user at the receiver side can see the quality measure instantaneously. While at the sender side, if necessary, video quality can be sent by the transport protocol from time to time (in RTCP protocol, it can be sent every 5 sec.). This means that the frequency update of the parameters and hence the quality evaluation can be done at any time the user wants at the receiver side, while at the sender side the user can have a feedback about the quality at least every 5 sec.

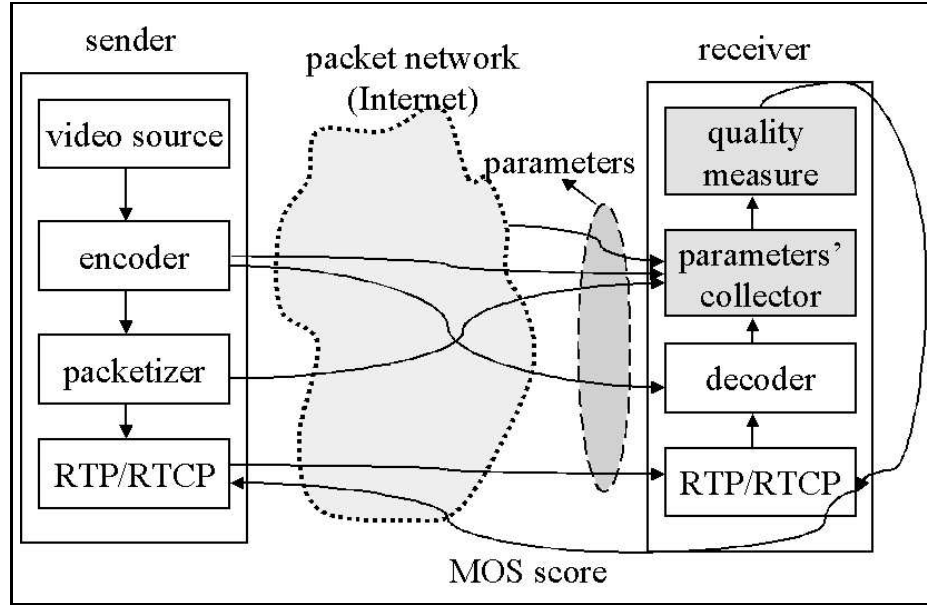


Figure 2: Operation mode for the tool in real-time video system.

## 4 Subjective Quality and MOS Calculations

To evaluate the quality of video systems (codec, telecommunication, television pictures, etc.), a subjective quality test is used. In this test, a group of human subjects is invited to judge the quality of the video sequence under the system conditions (distortions). There are several recommendations [12], [22] that specify strict conditions to be followed in order to carry out the subjective test. As mentioned in the introduction, the only reliable result is given by subjective quality tests. 1) Subjective quality methods The main subjective quality methods are Degradation Category Rating (DCR), Pair Comparison (PC) and Absolute Category Rating (ACR). For our case, we are using DCR method.

In the DCR subjective quality test, a pair of video sequences is presented to each observer, one after the other. They should see the first one, which is not distorted by any impairment, and then the second one, which is the original signal distorted by some configuration of the set of chosen quality-affecting parameters. Figure 3 shows the sequence and timing of presentations for this test. The time values come from the recommendation of the ITU-R [12].

As the observer is faced by two sequences, he/she is asked to assess the overall quality of the distorted sequence with respect to the non-distorted one (reference sequence). Figure 4 depicts

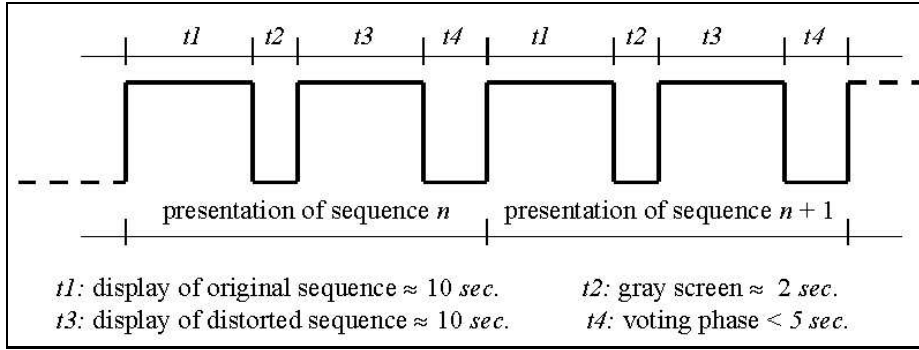


Figure 3: Presentation structure of video sequences to the set of human observers in a DCR subjective quality test experiment.

the ITU-R nine-grade scale. The observers should give the test presentation a grade from one to nine corresponding to their mental measure of the quality associated with it. It should be noted that there exist several quality scales. We chose this nine-grade one as a tradeoff between precision and dispersion of the subjective evaluations.

Following the ITU-R recommendations, overall subjective tests should be divided into multiple sessions and each session should not last more than 30 minutes. For every session, we should add several dummy sequences (about four or five). These sequences should not be taken into account in the calculation. Their aim is to be used as training samples for the observers to learn how to give meaningful rates.

#### 4.1 MOS Calculation

After performing any of the subjective methods, a range of integer values is given for each presentation. There will be variations in these distributions due to the differences in judgment between observers. For each of the test conditions, the mean score is given by:

$$\bar{u}_j = \frac{1}{N} \sum_{i=1}^N u_{ij}$$

where  $u_{ij}$  is the score of observer  $i$ , for test condition  $j$ , and  $N$  is the number of observers. For each test condition a *confidence interval* may be calculated. It is recommended to use the 95% confidence interval, given by:

$$[\bar{u}_j - \delta_j, \bar{u}_j + \delta_j]$$

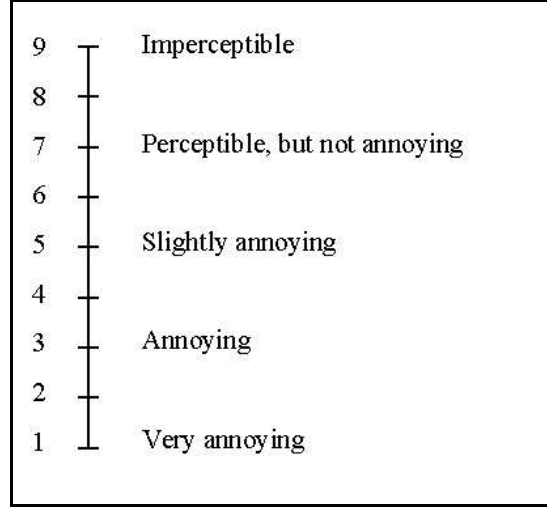


Figure 4: The ITU-R nine-grade standard quality scale.

where

$$\delta_j = 1.96 \frac{S_j}{\sqrt{N}}$$

and the estimated standard deviation,  $S_j$ , is given by

$$S_j = \sqrt{\sum_{i=1}^N \frac{(\bar{u}_j - u_{ij})^2}{N-1}}.$$

With a probability of 95%, the absolute value of the difference between the experimental mean score and the “true” mean score (for a very high number of observers) is then smaller than the 95% confidence interval.

## 4.2 Screening of the Observers

It must be ascertained whether this distribution of scores for test presentation is normal or not using the  $\beta_2$  test (by calculating the kurtosis coefficient of the function, i.e. the ratio of the fourth order moment to the square of the second order moment). The process can be expressed mathematically as follows. For each test condition calculate the mean, the standard deviation,

$S_j$ , and kurtosis coefficient,  $\beta_{2j}$ , where  $\beta_{2j}$  is given by:

$$\beta_{2j} = \frac{m_{4j}}{(m_{2j})^2}, \quad \text{with} \quad m_{xj} = \frac{\sum_{i=1}^N (u_{ij} - \bar{u}_j)^x}{N}.$$

For each observer  $i$ , find  $P_i$  and  $Q_i$  (initialized to 0) using [12]

if  $2 \leq \beta_{2j} \leq 4$ , then:

if  $u_{ij} \geq \bar{u}_j + 2S_j$  then  $P_i = P_i + 1$ ;

if  $u_{ij} \leq \bar{u}_j - 2S_j$  then  $Q_i = Q_i + 1$ ;

else

if  $u_{ij} \geq \bar{j} + \sqrt{20}S_j$  then  $P_i = P_i + 1$ ;

if  $u_{ij} \leq \bar{j} - \sqrt{20}S_j$  then  $Q_i = Q_i + 1$ .

If  $(P_i + Q_i)/J > 0.05$  and  $|(P_i - Q_i)/(P_i + Q_i)| < 0.3$  then reject observer  $i$ , where  $J$  is the number of test conditions.

## 5 Neural Networks

An ANN is a parallel distributed associative processor, comprised of multiple elements (neural models) highly interconnected [23]. Each neuron carries out two operations. First, we have an inner product of an input vector (carrying signals coming from other neurons) and a weight vector, where the weights represent the efficiencies associated with those connections and/or from external inputs. Second, the neuron has associated with, a nonlinear mapping between the inner product and a scalar, usually given by a non-decreasing continuous function (e.g., *sigmoid*, or *tanh*). When building ANN, an architecture and a learning algorithm must be specified. There are multiple architectures (e.g., multilayer feedforward networks, recurrent networks, bi-directional networks, etc.), as well as learning algorithms (e.g., Backpropagation, Kohonen's LVQ algorithm, Hopfield's algorithm, etc.).

As pattern classifiers, ANN work as information processing systems that search for a non-linear function mapping a set  $\mathcal{X}$  of input vectors (patterns) to a set  $\mathcal{Y}$  of output vectors (categories). This mapping is established by extracting the *experience* embedded in a set of examples (training set), following the learning algorithm. Thus, in developing an application with ANN, a set of  $N$  known examples must be collected, and represented in terms of patterns and categories (i.e., in pairs  $(X_n, Y_n)$ ,  $n = 1, 2, \dots, N$ , where  $X_n \in \mathcal{X}$  and  $Y_n \in \mathcal{Y}$ ). Then, an appropriate architecture should be defined. Last, a learning algorithm needs to be applied in order to build the mapping.

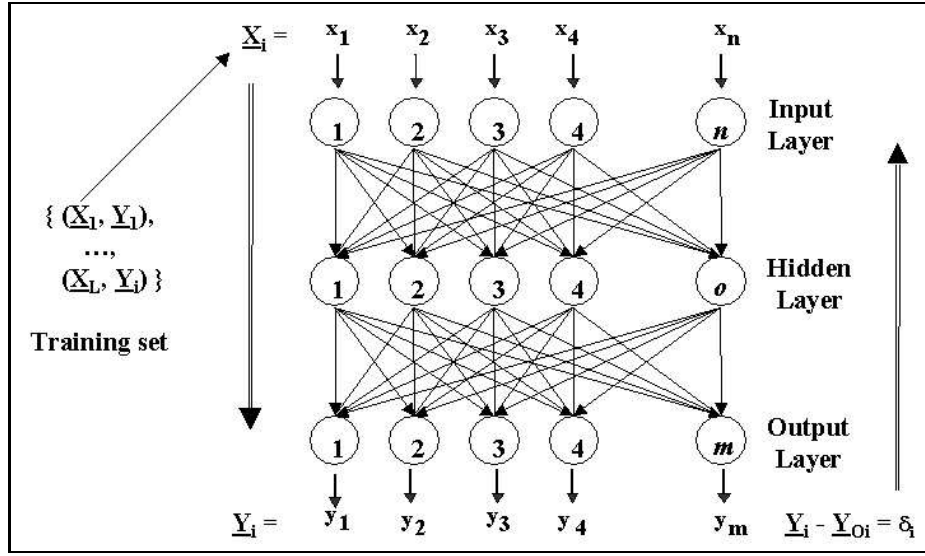


Figure 5: Architecture of a three-layer feedforward neural network.

Highly nonlinear mappings can be obtained using the backpropagation algorithm for learning and adaptation, and a three-layer feedforward neural network consisting of an input layer, a hidden layer and an output layer. In this architecture (see Figure 5), external inputs are the inputs for the neurons in the input layer. The scalar outputs from those neural elements in the input layer are the inputs for the neurons in the hidden layer. The scalar outputs in the hidden layer become, in turn, the inputs for the neurons in the output layer.

When applying the backpropagation algorithm, all weights initial values are given randomly. Then, for each pair  $(X_n, Y_n)$  in the database, the vector  $X_n$  (i.e., pattern) is placed as input for the input layer, and the process is carried out forward through the hidden layer, until the output layer response  $Y_{on}$  is generated. Afterwards, an error  $\delta$  is calculated by comparing vector  $Y_n$  (classification) with the output layer response,  $Y_{on}$ . If they differ (i.e., if a pattern is misclassified), the weight values are modified throughout the network accordingly to the generalized delta rule:

$$w_{ij}(t+1) = w_{ij}(t) + \delta \eta x_j,$$

where,  $w_{ij}$  is the weight for the connection that neuron  $i$ , in a given layer, receives from neuron  $j$  from the previous layer;  $x_j$  is the output of neuron  $j$  in that layer;  $\eta$  is a parameter representing the learning rate; and  $\delta$  is an error measure. In case of the output layer,  $\delta = ||Y_n - Y_{on}||$ , whereas in all hidden layers,  $\delta$  is an estimated error, based on the backpropagation of the errors



calculated for the output layer (for details refer to [23]). In this way, the backpropagation algorithm minimizes a global error associated with all pairs  $(X_n, Y_n)$ , where  $n = 1, 2, \dots, N$  in the database. The training process keeps on going until all patterns are correctly classified, or a pre-defined minimum error has been reached.

## 6 Results

In order to validate the applicability of our method, we chose to apply it to the assessment of subjective quality of real-time video transmission over IP networks.

### 6.1 Simulator Description

To generate the distorted video sequences, we used a tool that encodes a real-time video stream over an IP network into H263 format [13], simulates the packetization of the video stream, decodes the received stream and allows us to handle the simulated lost packets (for instance, for statistical purposes). The encoder can be parameterized, we can control the bit rate, the frame rate, the intra macro blocs refresh rate (i.e. encode the given macro bloc into intra mode rather than inter mode -this is done to make the stream resistant to losses [17]), image format (QCIF, CIF...), etc. The packetization process is in conformance with RFC 2429 [1].

We used a standard video sequence called *stefan* to test the performance of H26x and MPEG4 codecs. It contains 300 frames encoded into 30 frames per sec., and lasts for 10 secs. The encoded sequence's format is CIF (352 lines x 288 pixels). The maximum allowed packet length is 536 bytes, in order to avoid the fragmentation of packets between routers.

### 6.2 The Quality-affecting Parameters

We present here the quality-affecting parameters that we consider having the highest impact on quality:

- The bit rate (BR) in Kbps: this is rate of the actual encoder's output. It is chosen to take four values (256, 512, 768 and 1024 Kbps.). The effect of this parameter on quality is studied in [28].
- The number of frames per second (FR): the original video sequence is encoded at 30 frames per sec. This parameter takes one of 5 values (5, 10, 15 and 30 fps). This is done by the encoder by dropping frames uniformly. A complete study for the effect of this parameter on quality is given in [10].

- The ratio of the encoded intra macro-blocs to inter macro-blocs (G): this is done by the encoder, by changing the refresh rate of the intra macro-blocs in order to make the encoded sequence more or less sensitive to the packet loss [17]. This parameter takes values that vary between 0.053 and 0.4417 depending on the BR and the Fps for the given sequence. We selected five values for it.
- The packet loss rate (LR): the simulator can drop packets randomly and uniformly to satisfy a given percentage loss. This parameter takes five values (0, 1, 2, 4, and 8 %). It is admitted that a loss rate higher than 8 % will drastically reduce the video quality. In the networks where the LR is expected to be higher than this value, some kind of FEC [25] should be used to reduce the effect of losses. There are many studies analyzing the impact of this parameter on quality; see for example [4], [10], [11].
- The number of consecutively lost packets (CL): we chose to drop packets in bursts of 1 to 5 packets. These values come from real measurements that we performed before [20]. A study of the effect of this parameter upon the quality is, for instance, [11].

The delay and the delay variation are indirectly considered: they are included in the LR parameter. Indeed, it is known that if a dejittering mechanism with a strict playback buffer length is used, then all the packets arriving after a predefined threshold are considered as lost [8]. So, in this way, all delays and delay variations are mapped into loss.

If we choose to consider all the combinations of these parameters' values, we have to take into account  $4 \times 4 \times 5 \times 5 \times 5 = 2000$  different combinations. It is the role of the ANN to interpolate the quality scores for the missed parts of this potentially large input space. We chose to give default values and to compose different combinations by changing only two parameters at a time. This led to 94 combinations.

### 6.3 Subjective Quality Test and MOS Experiment

The subjective quality test is with conformance to the method Degradation Category Rating (DCR), with a quality scale consisting of 9 points, see Section D. We divided the test into two sessions, and added 5 distorted sequences to the first session and 4 to the second session. These nine sequences will not be considered in the MOS calculation as their aim is to be used as a training phase for the human subjects. At the same time, they are used to verify how much reliable is the person carrying out the test, as they are replicated from the real 94 samples.

We invited 20 persons to perform the subjective tests. After that, a prescreening of the results was performed; as a consequence, we discarded the notes of two subjects. The 95% confidence intervals after and before removing these two subjects are shown in Figure 6. As it is clear from

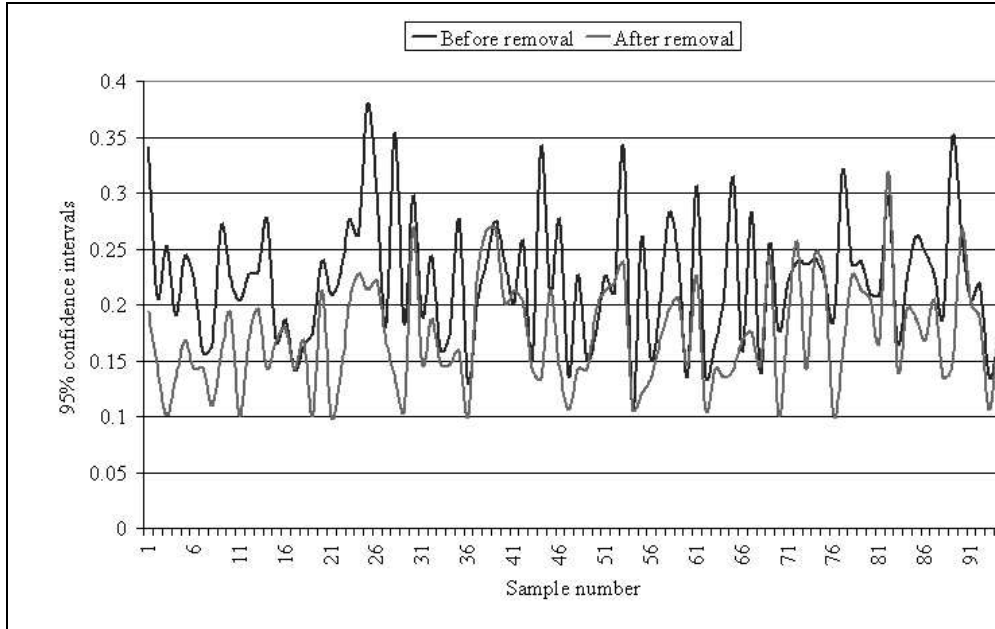


Figure 6: The 95% confidence intervals before and after removing the rates of two unreliable subjects. We can see how the size of the confidence interval decreases after the removal of the rates of these two persons.

the figure, the removal of the rates of these two persons significantly increased the precision of the overall MOS scores.

#### 6.4 Training and Testing the ANN

The number of input neurons in the input layer of our ANN is equal to the number of selected parameters (five). There is only one output of the ANN, the MOS score.

The number of hidden neurons in the hidden layer is variable as it depends on the complexity of the problem (inputs, outputs, training set, and the global precision needed). There is a trade-off between the number of hidden neurons and the ability of the ANN to generalize (i.e. to produce good results for inputs not seen during the training phase). Hence, the number of hidden neurons should not be too large. In our case we chose eight neurons.

After carrying out the MOS experiment for the 94 samples, we divided our database into two parts: one to train the ANN containing 80 samples, and the other to test the performance

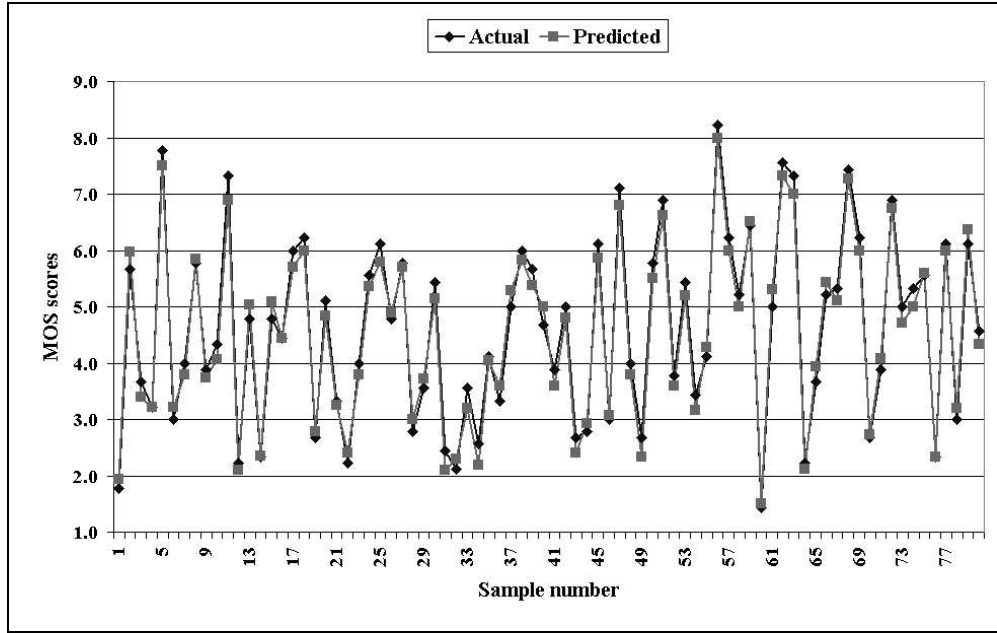


Figure 7: Actual and Predicted MOS scores for the training database. The neural network learned the way by which the group of human subjects rated video quality for the video samples with very good precision.

of the ANN to work in a dynamic environment, containing 14 samples. After training the ANN using the first database and comparing the training set against the values predicted by the ANN, we got a correlation factor = 0.9915, and average absolute error = 0.2084; that is, the neural network model fits quite well the way in which humans rated the video quality. The result is shown in Figure 7.

### 6.5 How Well does the ANN Perform?

In order to address the question “How well does the ANN perform?”, the ANN was applied to the testing set (which contains samples that never have been seen during the training process). The results were correlation coefficient = 0.9907 and average error = 0.253. Once again the performance of the ANN was excellent, as can be observed in Figure 8.

>From Figure 7 and Figure 8, it can be observed that the video quality scores generated by the ANN fits quite nicely the nonlinear model built by the subjects participating in the MOS

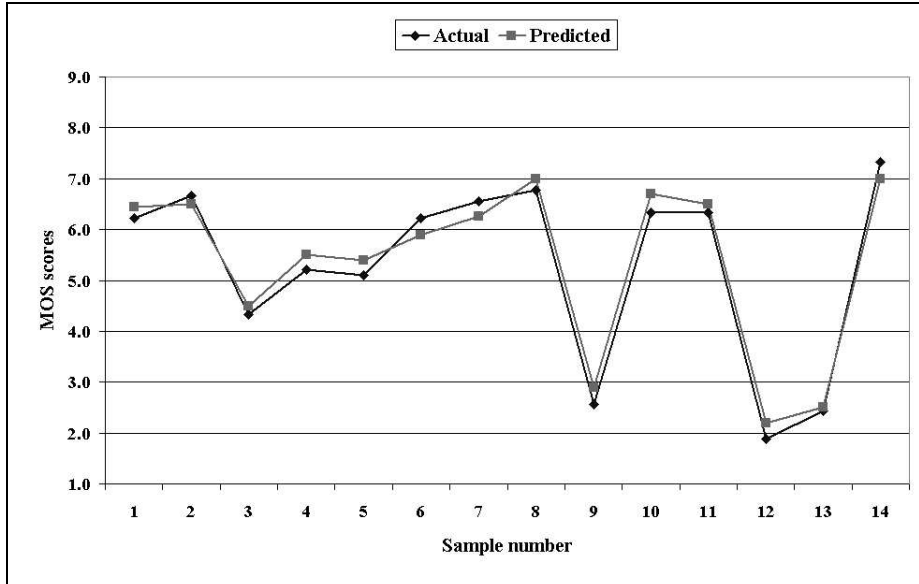


Figure 8: Actual and Predicted MOS scores for the testing database. The trained neural network is able to evaluate video quality for new samples that never have been seen during the training, with good correlation with the results obtained by human subjects

experiment. Also, by looking at Figure 8, it can be established that learning algorithms give neural networks the advantage of high adaptability, which allows them to self optimize their performance when functioning under a dynamical environment (that is, reacting to inputs never seen during the training phase).

## 7 Conclusion and Future Directions

In this paper, it has been described how ANN can be used to create a nonlinear mapping between non subjective audio signals measures (i.e., packet loss rate, loss distribution, bit rate, frame rate, encoded frame type, etc.), and a subjective (i.e., MOS) measure of video quality. This mapping mimics the way in which human subjects perceive video quality at a destination point in a communication network.

We have validated our approach by building the INN to assess in real time the video quality transmitted over the Internet, taking into account the previously mentioned parameters. We have shown that the ANN performs quite well in measuring video quality in real time.

As the video quality is affected by many parameters, one of the future directions in our research is to build a robust database, by conducting a series of MOS experiments taking into account different combinations of these parameters. The ANN approach allows also to identify the importance of network parameters in distorting video signals. Thus, using a tool that effectively measures video quality and identifies the nature of current distortions, better solutions to other problems would be developed, e.g., adaptive error correction schemes to dynamically compensate video distortion based on the current network situation, identification of the best trade-off between redundant information and bandwidth requirements to improve QoS, etc.

Our method is applied to packet networks, but can be also applied to other network technologies. It will be interesting to further investigate these possibilities in the future.

## References

- [1] 2429, R. RTP payload format for the 1998 version of ITU-T rec. H.263 video (H.263+). In *IETF* (Oct. 1998).
- [2] ABDELBAKI, H., GELENBE, E., AND EL-KHAMY, S. Random neural network decoder for error correcting codes. In *International Joint Conference on Neural Networks – IJCNN '99* (Vol. 5 , 1999, pp. 3241 -3245).
- [3] ANTILA, I. Transferring real-time video on the internet. In *Proceedings of the HUT Internetworking Seminar* (May 1997).
- [4] BOYCE, J., AND GAGLIANELLO, R. Packet loss effects on MPEG video sent over the public internet. In *Proceedings of ACM Multimedia'98* (1998).
- [5] CARAMMA, M., LANCINI, R., AND MARCONI, M. Subjective quality evaluation of video sequences by using motion information. In *Proceedings of International Conference on Image Processing: ICIP'99* (Kobe, Japan, 24-28 October 1999).
- [6] CLAYPOOL, M., AND TANNER, J. The effects of jitter on the perceptual quality of video. In *Proceedings of ACM Multimedia Conference* (1999).
- [7] CRAMER, C., GELENBE, E., AND GELENBE, P. Image and video compression. *IEEE Potentials* (February/March 1998).
- [8] DE VLEESCHAUWER, D., JANSSEN, J., AND PETIT, G. H. Delay bounds for low bit rate voice transport over IP networks. In *Proceedings of the SPIE Conference on Performance and Control of Network Systems III* (Vol. 3841, pp. 40-48, Boston (MA), 20-21 Sept. 1999).

- [9] GELENBE, E., BAKIRCIOGLU, H., AND KOC AK, T. Image processing with the random neural network (RNN). In *13th International Conference on Digital Signal Processing Proceedings - DSP'97* (1997, pp. 243 -248).
- [10] GHINEA, G., AND THOMAS, J. QoS impact on user perception and understanding of multimedia video clips. In *Proceedings of ACM Multimedia 98* (1998).
- [11] HANDS, D., AND WILKINS, M. A study of the impact of network loss and burst size on video streaming quality and acceptability. In *Interactive Distributed Multimedia Systems and Telecommunication Services Workshop* (Germany, 1999).
- [12] ITU-R RECOMMENDATION BT.500-10. Methodology for the subjective assessment of the quality of television pictures. <http://www.itu.int/>.
- [13] ITU-T RECOMMENDATION H.263. Video coding for low bit rate communication. <http://www.itu.int/>.
- [14] JONES, C., AND ATKINSON, D. Development of opinion-based audiovisual quality models for desktop video-teleconferencing. In *International Workshop on Quality of Service* (1998).
- [15] LAVINGTON, S., DEWHURST, N., , AND GHANBARI, M. The performance of layered video over an IP network. In *submitted to Packet Video* (2000).
- [16] LAVINGTON, S., HAGRAS, H., AND DEWHURST, N. Using a MLP to predict packet loss during real-time video transmission. In *Univ. Of Essex, UK, Internal Report CSM329* (July 1999).
- [17] LE LEANNEC, F., AND GUILLEMOT, C. Packet loss resilient H.263+ compliant video coding. In *Proceedings of International Conference on Image Processing* (Sept. 2000).
- [18] MITCHELL, M. *An Introduction to Genetic Algorithms*. MIT Press, Cambridge, MA, 1996.
- [19] MOHAMED, S., CERVANTES, F., AND AFIFI, H. Integrating networks measurements and speech quality subjective scores for control purposes. In *Proceedings of IEEE INFOCOM'01* (April 22 - 26, 2001, Alaska, USA).
- [20] MOHAMED, S., CERVANTES, F., AND AFIFI, H. Audio quality assessment in packet networks: an inter-subjective neural network model. In *Proceedings of the IEEE 15th International Conference on Information Networks: ICOIN-15* (Japan, Jan. 2001).
- [21] PESSOA, A., FALCAO, A., NISHIHARA, R., SILVA, A., AND LOTUFO, R. Video quality assessment using objective parameters based on image segmentation. *SMPTE Journal* (Dec. 1999).

- 
- [22] REC. ITU-T P.910. Subjective video quality assessment methods for multimedia applications. <http://www.itu.int/>.
  - [23] RUMELHART, D., HINTON, G., AND WILLIAMS, R. *Learning internal representations by error propagation*. Parallel Distributed Processing, vol. 1. Cambridge, Massachusetts, MIT Press, 1986.
  - [24] TAN, D., AND A., Z. Real-time internet video using error resilient scalable compression and TCP-friendly transport protocol. *IEEE Transactions on Multimedia* (May 1999).
  - [25] TAN, K., AND GHANBARI, M. A multi-metric objective picture-quality measurement model for MPEG video. *IEEE Transactions on Circuits and Systems for Video Technology* (vol. 10, no. 7, Oct. 2000, pp. 1208 -1213).
  - [26] WATSON, A., AND SASSE, M. Evaluating audio and video quality in low-cost multimedia conferencing systems. *Interacting with Computers* 8, 3 (1996), 255–275.
  - [27] WATSON, A., AND SASSE, M. Measuring perceived quality of speech and video in multimedia conferencing applications. *Proceedings of ACM Multimed'98* (pp. 55-60, 12-16 Sept., 1998).
  - [28] WU, H., FERGUSON, T., AND QIU, B. Digital video quality evaluation using quantitative quality metrics. In *Proceedings of the 4th International Conference on Signal Processing* (1998, Oct.).
  - [29] WU, H., LAMBRECHT, C., YUEN, M., AND QIU, B. Quantitative quality and impairment metrics for digitally coded images and image sequences. In *Proceedings of Australian Telecommunication Networks & Applications Conference* (Dec. 1996).





---

Unité de recherche INRIA Lorraine, Technopôle de Nancy-Brabois, Campus scientifique,  
615 rue du Jardin Botanique, BP 101, 54600 VILLERS LÈS NANCY  
Unité de recherche INRIA Rennes, Irisa, Campus universitaire de Beaulieu, 35042 RENNES Cedex  
Unité de recherche INRIA Rhône-Alpes, 655, avenue de l'Europe, 38330 MONTBONNOT ST MARTIN  
Unité de recherche INRIA Rocquencourt, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex  
Unité de recherche INRIA Sophia-Antipolis, 2004 route des Lucioles, BP 93, 06902 SOPHIA-ANTIPOLIS Cedex

---

Éditeur  
INRIA, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex (France)  
<http://www.inria.fr>  
ISSN 0249-6399